



caBIG

*cancer Biomedical
Informatics Grid*



Data and Vocabulary Standards

October 26, 2004

Kathleen Gundry (SAIC)

Tommie Curtis (SAIC)

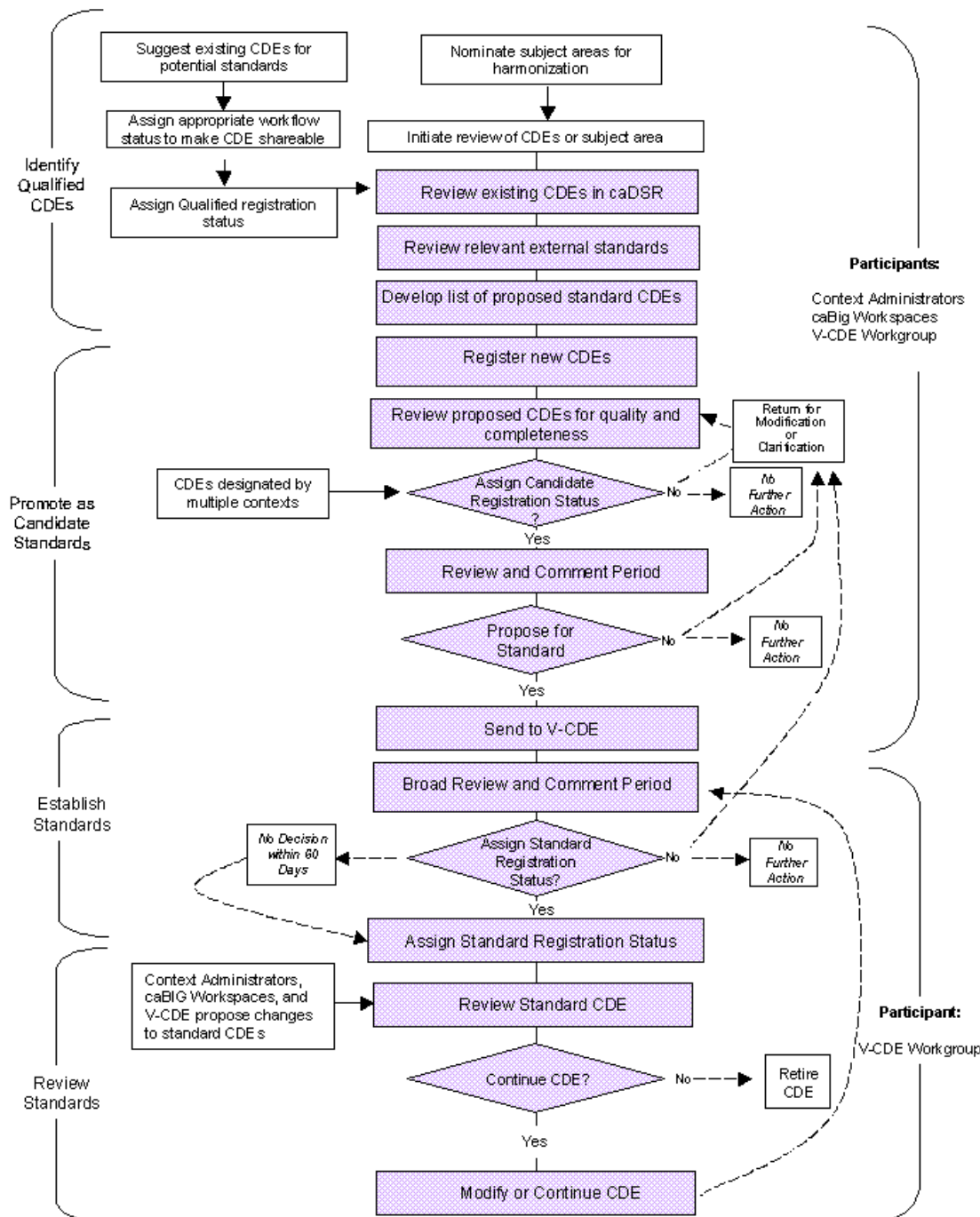
Contractors for NCICB

Agenda

2

- ▶ Draft Governance for Review and Approval of caBIG Data and Vocabulary Standards
- ▶ Review Potential Candidate Data and Vocabulary Standards
- ▶ Next Steps for V-CDE

Draft caBIG Data Standard Governance Process



NCI External Standards Review

4

- ▶ SAIC compiled a review of data standards relevant to the caBIG.
- ▶ The collection includes both vocabulary/coding standards as well as exchange format standards.
- ▶ The report makes recommendations for caBIG use of the standards. The table categories them as:
 - A - Recommended for caBIG use.
 - B - Recommended for further consideration by caBIG.
 - C - Standard is not recommended at this time.

<ftp://ftp1.nci.nih.gov/pub/cacore/ExternalStds/>

Types of Standards

5

- ▶ Content – controlled vocabularies, ontologies, value lists.
- ▶ Information Collection and Format – Representation of Date and Time, PEDro, MIAME.
- ▶ Exchange/Transaction – XML, PPI.

Overview: Issues Associated with Standards

6

- ▶ Fulfilling a need/niche, timing.
- ▶ Community participation – create/use.
- ▶ Ease of use, availability of software tools.
- ▶ Accommodate change.
- ▶ New standards are being created.

Common Demographic/Information Processing and Code Sets

7

Standard Type	Standard Name	Standard Content	Recommendation (EVS Entry)
Common Demographic / Information Processing and Code Sets			
Address	FIPS 5-2, Codes for Identification of the States, District of Columbia, and Outlying Areas of the United States, and Associated Areas	State Names	A
	FIPS 10-4, Countries, Dependencies, Areas of Special Sovereignty, and their Principal Administrative Divisions	Country Names	A
	ISO 3166-1, Country Codes	Country Codes	A
	ISO 11180:1993, Postal Addressing	Address Format	A
	Universal Postal Union	Address Format, State Codes, Country Codes	A
	U.S. Postal Service Postal Addressing Standards	Address Format, State Codes, Street Suffixes, Secondary Unit Designators	B

Common Demographic/Information Processing and Code Sets

8

Standard Type	Standard Name	Standard Content	Recommendation (EVS Entry)
Common Demographic / Information Processing and Code Sets			
Language	ISO 639, Codes for representation of language	Language Codes	A
Race and Ethnicity	Office of Management and Budget Directive 15, Standards for the Classification of Federal Data on Race and Ethnicity	Race Identification, Ethnicity Identification	A
Occupation Classification	Bureau of Labor Statistics, Standard Occupational Classification System	Job Activity Classification	A
Vital Statistics	Centers for Disease Control (CDC) National Center for Health Statistics	Birth and death records, Medical records, Interview surveys, Physical exams, Laboratory testing, Marriages and divorces, Fetal death	A
Measurement	HL7 codes for Units, Versions 2.X + (derived from the ISO 2955-83 standard [withdrawn by ISO in 2001] and ANSI X3.50)	Common units of measure, such as Celsius or mg/ml, intended to be combined with a numeric value to accurately express a result	A
	ISO 31, Quantities and units	Individual standards dealing with quantities in space and time, periodic phenomena, mechanics, heat, electricity and magnetism, electromagnetic radiation, chemistry, molecular physics, nuclear physics	A
Information Processing	FIPS 4-2, Representation of Calendar Date for Information Interchange	Means of representing calendar date to facilitate interchange of data among information systems	A
	ISO 8601, Numeric representation of dates and times	Formats for date and time	A

Health-related Vocabulary/Coding Standards

9

Standard Type	Standard Name	Standard Content	Recommendation (EVS Entry)
Health-related Vocabulary / Coding Standards			
Health Thesaurus	National Cancer Institute (NCI) Thesaurus	NCI reference terminology and description-logic ontology, providing comprehensive classification and characterization of types of cancer as well as cancer-related diseases, disorders, findings, abnormalities (cellular, molecular, and cytogenetic), gross anatomy, microanatomy, biological processes, genes, gene products, chemicals/drugs, combination therapies, mouse and other experimental models, and other topics	A
Basic Biology	Biological Pathways Exchange (BioPax)	Ontology for pathway information	B
	International Union of Biochemistry and Molecular Biology (IUBMB) and the International Union of Pure and Applied Chemistry (IUPAC)	Controlled vocabulary for nomenclature for biochemistry and molecular biology	A

Health-related Vocabulary/Coding Standards

10

Standard Type	Standard Name	Standard Content	Recommendation (EVS Entry)
Health-related Vocabulary / Coding Standards			
Clinical	Common Terminology Criteria for Adverse Events v 3.0 (CTCAE)	Descriptive terminology for Adverse Event reporting	A
	Current Procedural Terminology (CPT) 4	Coding for evaluation and management, anesthesia, surgery, radiology, pathology and laboratory, medicine	B
	Healthcare Common Procedure Coding System (HCPCS)	Healthcare procedures, equipment, and supplies (Level 1) - used for Medicare billing Classification (national level) of physician and non -physician patient care services (Level 2)	B (HCPCS04)
	International Classification of Diseases for Oncology (ICD-O-3)	Coding for diagnoses of neoplasms - both topography and morphology - includes tumor location, cell type, tumor type, aggressiveness grade	A (ICS03, 3 rd Ed, 2000)
	International Classification of Diseases, Clinical Modification (ICD-9-CM)	Classifies diseases, conditions, symptoms, complaints/problems by diagnosis; supplementary classifications include health status, external causes of injury and poisoning, morphology of neoplasms, glossary of mental disorders, drug list numbers, industrial accidents and surgical, diagnostic, and therapeutic procedures	B (ICD9CM_2004)

Health-related Vocabulary/Coding Standards

11

Standard Type	Standard Name	Standard Content	Recommendation (EVS Entry)
Health-related Vocabulary / Coding Standards			
Clinical	International Statistical Classification of Diseases and Related Health Problems (ICD-10)	Collection, processing, classification, and presentation of mortality statistics	A (ICS-10-1998)
	Logical Observation Identifiers Names and Codes (LOINC)	Standard test names and codes, descriptive elements for other healthcare areas	A (LOINC No version provided)
	Medical Dictionary for Regulatory Activities (MedDRA)	Signs, symptoms, diseases, diagnoses, therapeutic indications, names and qualitative results, surgical and medical procedures, medical/social/family history, adverse event reporting	A (MED-60 Version 6 March 2003)
	Systematized Nomenclature of Human and Veterinary Medicine (SNOMED)	Findings/conclusions/assessments, procedures, body structures, function, organisms, substances, physical agents, occupations, social context/demographics, specimens, and other concepts	A (SNOWMEDCT_2004_01_31)

Health-related Vocabulary/Coding Standards

12

Standard Type	Standard Name	Standard Content	Recommendation (EVS Entry)
Health-related Vocabulary/Coding Standards			
Genomics	Gene Ontology (GO)	Structured, controlled vocabularies describing gene products used for gene annotations	A
	HUGO Gene Nomenclature Committee (HGNC)	Controlled vocabulary of gene names and symbols for human genes	A
	Mammalian Phenotype Ontology (MP)	Standard vocabulary to describe phenotype data	B
	The Microarray Gene Expression Data (MGED) Society	Comprised of MIAME, MAGE, and the MAGE ontology, a suite of standards for microarray users and developers including an object model, document exchange format, toolkit, and ontology	A
	Mouse Anatomy (adult – MA, and development – EMAP)	Ontologies used to annotate gene products	A
	Taxonomy	National Center for Biotechnology Information (NCBI) taxonomy of organism names represented in genetic databases	A
Drug Identification	National Drug File Reference Terminology (NDF-RT)	Drug classes, active ingredients (chemical structure), mechanics of action, physiologic effect, pharmacokinetics, therapeutic intent, commercial/clinical drug identification	A (VANDF03)
	RxNorm Clinical Drug Vocabulary	Ingredients, drug components, drug formulations, drug strength representation, drug name synonyms, dosage forms	A (RxNORM_04AA)

Health-related Transaction Standards and Models

13

Standard Type	Standard Name	Standard Content	Recommendation (EVS Entry)
Health-related Transaction Standards and Models			
Imaging	Digital Imaging and Communications in Medicine	Standard method for the transmission of medical images and their associated information	A
Basic Biology	Systems Biology Markup Language (SBML)	XML exchange format for exchange of biochemical network models	A
	CellML	XML-based language for describing and exchanging models of cellular and subcellular processes	A
Clinical	Health Level Seven (HL7)	Patient tracking, scheduling, orders, results, clinical observations, billing, medical records, patient referral, patient care	A (HL7_1998-2002 Vocabulary)
	Clinical Data Interchange Standards Consortium (CDISC)	Clinical trials data - general data (study name, protocol name, measurement units), study metadata (code lists), administrative data, reference data (lab normal ranges), clinical data	A
	North American Association of Central Cancer Registries, Inc. (NAACCRz, Inc.)	Demographic, tumor and staging, treatment and follow-up	A
Genomics / Proteomics	Macromolecular Structure (Mms)	Specification for a data model and interface for exchange of macromolecular structure information	B
	PEDRo	Data model implemented in SQL and XML to support proteomics research	B
	Protein-Protein Interaction (PPI)	Data exchange format designed to bridge different formats of protein interaction databases	B
	Tissue Microarray (TMA)	Data exchange specification for tissue microarray data	A

Standards Review Bodies

14

Standard Type	Standard Name	Standard Content
Standards Review Bodies		
	Centers for Disease Control Public Health Information Network (PHIN)	Vocabulary and messaging standards; standards for data display and entry; standards for data transmission and management; implementation of applications and databases to support the adopted data standards
	Consolidated Health Informatics Initiative (CHI)	Portfolio of existing clinical vocabularies and messaging standards enabling federal agencies to build interoperable federal health data systems
	Food and Drug Administration, Center for Drug Evaluation and Research (CDER)	Compilation of standardized nomenclature monographs that have been reviewed and approved by the CDER Nomenclature Standards Committee (NSC)
	National Council on Vital Health Statistics (NCVHS)	Advises the government on recommended standards for adoption in the health care sector

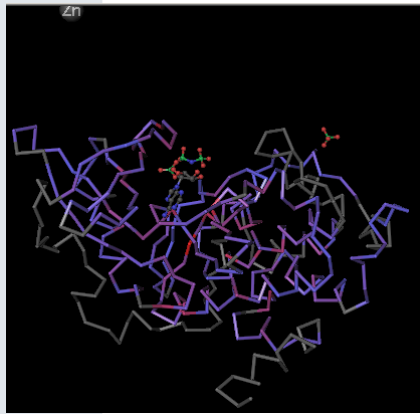
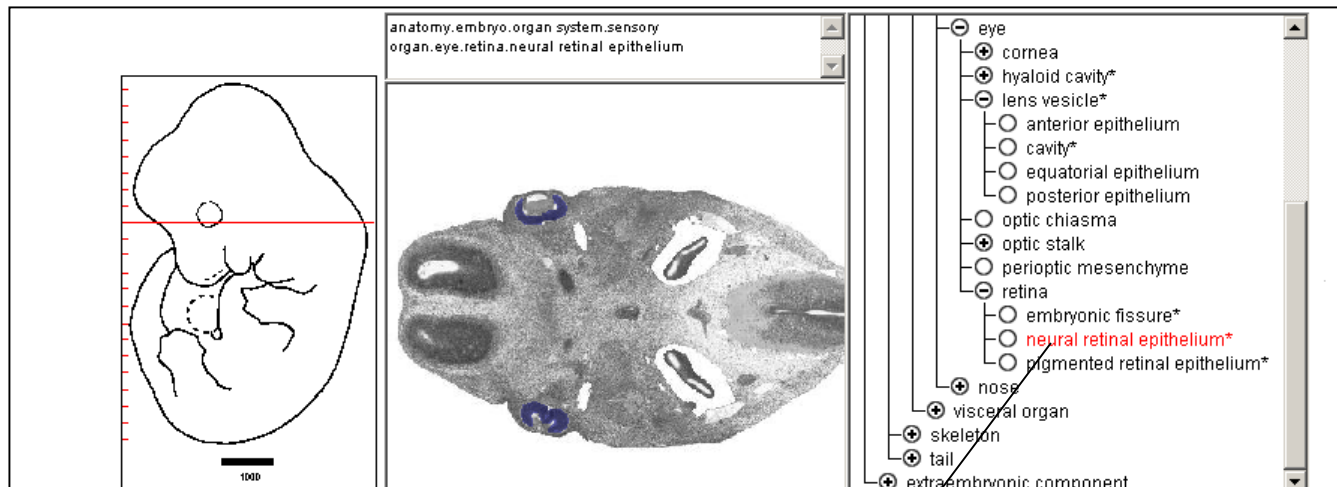
Overview of Reviewed External Biology Standards

15

<i>System</i>	TAX MP	TMA	<i>Tissue Array</i>
<i>Organic</i>	MA/EMAP MPATH		<i>Metabolomic</i>
<i>Cellular</i>	SBML/CellML/BioPAX GO	PEDRo PSI-MS PPI <i>PSI-OM/-ML</i> <i>/-Ont/MIAPE</i>	<i>Proteomic</i>
<i>Molecular</i>	PPI/Mms HGNC BSA/GM/SO IUPAC-IUBMB	MAGE MGED Ont MIAME Tox-MIAME	<i>Gene Exp.</i>

Example of Data Integration: MA & GXD

16



NCBI Conserved Domain Database

CD: [pfam00069.10.pkinase](#) PSSM-Id: 22891

Description: Protein kinase domain.

Taxa: [root](#)

Status: Alignment from source

Aligned: 68 rows

Proteins: [Click here for CDART summary of Proteins containing pfam00069](#)

View 3D Structure with [Cn3D](#) using [Virtual Bonds](#) (To display structure, download [Cn3D](#))

View Alignment as [Hypertext](#) width 60 color at 2.0 bits

Subset Rows up to 10 of the most diverse members

consensus 1 YELGEKLGSGSPGKVKQKHK-----GTGEIVAVKILK--KRSIKEkk-----rFLR 45

1JWH_A 39 YQLVRKLGKRGKYSEVFEAINI-----TNNEKVVVKILK--PVKKKK-----IKR 80

Gene Expression Data						
Query Results -- Summary						
43 matching assay results displayed						
Gene	Assay Type	Result Details	Mutation	Age	Structure	Detected?
Acvrb	RNA in situ	MGI3043600		E12.5	TS20: neural retinal epithelium	yes
	TA in situ	MGI3043600		E12.5	TS20: neural retinal epithelium	yes
	TA in situ	MGI2135922		E12.5	TS20: neural retinal epithelium	no
	TA in situ	MGI2181802		E12.5	TS20: neural retinal epithelium	yes
	TA in situ	MGI2181814		E12.5	TS20: neural retinal epithelium	yes
	TA in situ	MGI1342268		E12.5	TS20: neural retinal epithelium	yes

Next Steps for Data Standard Process

17

- ▶ Mentoring and outreach to domain workspaces to support identification of data and vocabulary standards needed for activities within a workspace.
- ▶ Identify data and vocabulary standards to meet the overarching needs across caBIG.
- ▶ Develop V-CDE process for review, approval, and maintenance of data and vocabulary standards for the caBIG context.
- ▶ How will data standards compliance and conformance be measured.

Examples of Data Standard Under Development by NCI Context Administrators

18

- ▶ Exchange Format for Date and Time Information
- ▶ Specification of Race and Ethnicity
- ▶ Specification of Gender
- ▶ Specification of Age
- ▶ Representation of Name, Organization, and Address Information
- ▶ Patient Performance Status Standards
- ▶ Protocol Identification Metadata
- ▶ Reporting Laboratory Results
- ▶ Serious Adverse Events Reporting

V-CDE Data Standards Review Process Development – Step 1

19

- ▶ V-CDE Receives Candidate Data Standards and Vocabularies
 - What is the format for submission?
 - What must be included in the package?
 - Who will receive the information?

V-CDE Data Standards Review Process Development – Step 2

20

- ▶ Implement Broad Review and Comment Period
 - What is the initial level of review by V-CDE?
 - Who will be included in the broader review?
 - How will the broader review be accomplished?
 - How will feedback/comments be addressed?
 - How long is the review period?
 - What are the criteria for progressing to Standard?

V-CDE Data Standards Review Process Development – Step 3

21

- ▶ On-going Maintenance of Data
 - Will there be periodic reviews of data and vocabulary standards?
 - Who will review?
 - How will update requests be handled?
 - How will change notification be provided?

Questions and Answers

22